

Vendi-RAG Performance Scaling with Document Corpus Size on TriviaQA

Assignee Research

May 30, 2026

Abstract

This report synthesises findings from 3 peer-reviewed papers addressing the following research question: How does Vendi-RAG's performance scale with increasing document corpus size in terms of EM score and latency on the TriviaQA benchmark compared to traditional RAG. The rapid evolution of natural language processing technologies has significantly enhanced the capabilities of generative models, yet challenges remain in maintaining the accuracy and relevance of information over time. The novel concept of integrating a hierarchical Retrieval. 4 claims were extracted from source literature; 4 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 8.2/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Improving Real-Time Knowledge Retrieval in Large Language Models with a DNS-Style Hierarchical Query RAG. Research question: How does Vendi-RAG's performance scale with increasing document corpus size in terms of EM score and latency on the TriviaQA benchmark compared to traditional RAG?.

2 Methodology

Systematic literature search across multiple databases yielded 3 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.2/10.

3 Results

3 papers retrieved. 4 claims extracted; 4 independently verified. Quality review score: 8.2/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
The hierarchical RAG system, designed with multilevel caching and dynamic query routing mechanisms, optimizes informatio	✓	0.37
Comprehensive evaluations demonstrated that the RAG-enhanced Llama outperformed the baseline model across various metric	✓	0.32
The system's adaptability to varying data loads and its efficient management of large-scale databases further underscore	✓	0.26
Comparisons with traditional retrieval methods revealed the distinct advantages of the hierarchical RAG system, particul	✓	0.31

References

- <https://doi.org/10.32628/ijrsrset242439>
- <https://doi.org/10.36227/techrxiv.171838950.05945972/v1>
- <https://doi.org/10.3390/sym16111470>